

$$x \sim p(x) = N(0, R)$$

covariance matrix  
 outer product  
 scalar

R is square & invertible

$$R = E[xx^T] = E\left[\begin{array}{|c|} \hline x \\ \hline \end{array} \begin{array}{|c|} \hline x^T \\ \hline \end{array}\right] \rightarrow R_{ij} = E[x_i x_j]$$

$$a^T R a = a^T E[xx^T] a = E[(a^T x)(x^T a)] = E[z^2] > 0$$

↳ R is positive definite

$$R \hat{=} \hat{R} = \frac{1}{N} \sum_{i=1}^N y_i y_i^T$$

## Multivariate Gaussian Distribution

Used a lot, in M.L. and elsewhere

- Let  $x$  be a zero-mean random variable on  $\mathbb{R}^p$

General Form  $p(\vec{x}) = \frac{1}{(2\pi)^{p/2}} |R|^{-1/2} \exp \left\{ -\frac{1}{2} \underbrace{\vec{x}^T R^{-1} \vec{x}}_{\text{quadratic form}} \right\}$

as exponent gets larger,  $p(\vec{x})$  gets smaller  
 $\vec{x} \in \mathbb{R}^p$        $R^{-1}$

where  $R$  is the  $p \times p$  covariance matrix.

$$\dots \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^T R^{-1} (\vec{x} - \vec{\mu}) \right\}$$

- The matrix  $R$  is a positive definite symmetric matrix, then

$$R = E \Lambda E^T$$

↳ because  
 $R_{ij} = E[x_i x_j] = E[x_j x_i] = R_{ji}$

where  $E = [e_1, \dots, e_p]$  is an orthonormal matrix of eigenvectors, and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  is a diagonal matrix of eigenvalues. ( $\lambda$  is positive)

- Therefore, we have  $E^T E = I$ .

$$E = [e_1, \dots, e_p]$$

$$e_i^T e_j = \delta(i-j) \rightarrow \text{orthonormal matrix}$$

$$R > 0 \text{ i.i.f. } \lambda_i > 0 \forall i$$

$$\begin{aligned}
 Re_i &= E \Lambda E^T e_i = E \Lambda \begin{bmatrix} e_1^T \\ \vdots \\ e_p^T \end{bmatrix} e_i \\
 &= E \Lambda \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} = E \begin{bmatrix} 0 \\ \vdots \\ \lambda_i \\ 0 \end{bmatrix} = \lambda_i e_i
 \end{aligned}$$

$$e^{j\omega n} \rightarrow \boxed{LTI} \rightarrow \lambda_n e^{j\omega n}$$

$$H(\omega) = \lambda_\omega$$

sine waves are eigenfunctions for LTI systems  
(eigenvectors)

Re-Watch:

- integral of  $p(\vec{x})$  has to be 1, that's why there is a coefficient in front of the  $\exp\{\cdot\}$

-  $R$  is symmetric & positive definite  
 $R = R^+$ ,  $\forall x \in \mathbb{R}^p, x^T R x > 0$

$$E^+ E = I : \begin{bmatrix} e_1^+ \\ \vdots \\ e_p^+ \end{bmatrix} [e_1, \dots, e_p] = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$$

$$e_i^T e_j = \begin{cases} 1, & i=j \\ 0, & i \neq j \end{cases} = \delta(i-j)$$

$R = E \Lambda E^+$  - works because  $R$  is symmetric

-  $E$  &  $\Lambda$  are real-valued

Idea:  $R e_i = \lambda_i e_i$

$E$ -eigentransform

$\Lambda$ -eigenvalues, all elements are greater than zero ( $R$  positive definite)

$$R^{-1} = (E \Lambda E^+)^{-1}$$

can invert R just by inverting the eigenvalues

$$= E \Lambda^{-1} E^+$$

## Contour for the Energy Function

- Intuitively, the energy function (square of the Mahalanobis distance)

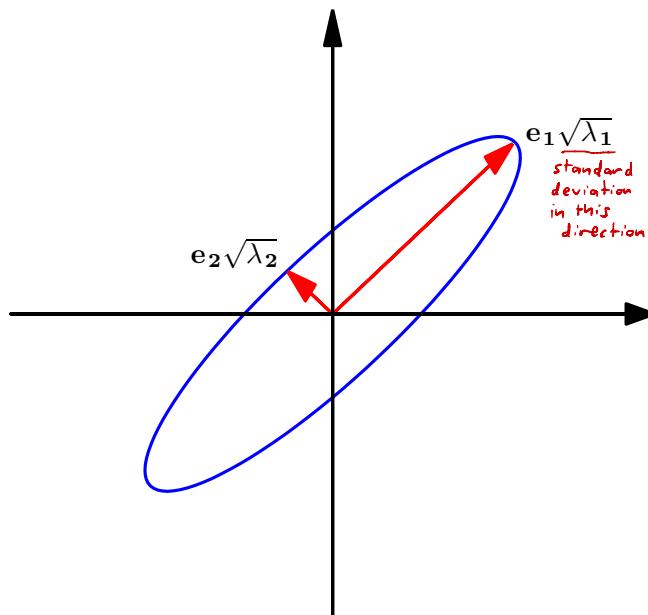
$$f(\mathbf{x}) = \mathbf{x}^t R^{-1} \mathbf{x}$$

$f(0) = 0$  (minimum)  
quadratic form

has contour plots shown here.

let  $\lambda_1 \geq \lambda_2$

Contour for the energy function



Scalar Case:  $\rho(x) = \frac{1}{2} \exp\left\{-\frac{1}{2\sigma^2} x^2\right\}$

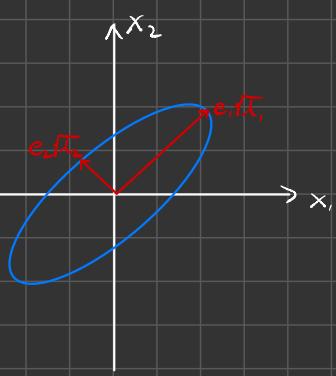
if eigenvalues negative, standard deviation is negative,  
which doesn't make sense

$$f(\vec{x}) = \vec{x}^T R^{-1} \vec{x}$$

$$f(\vec{0}) = 0$$

$$f(\vec{x}) = -\log(\rho(\vec{x})) + c$$

as  $\rho(\vec{x})$  goes down,  $f(\vec{x})$  increases



$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$$

$$\lambda_i \geq \lambda_{i+1}$$

- eigenvectors  $e_i$  are principal axes
- eigenvalues specify the length to travel along the principal axes

- if  $R$  diagonal, components are uncorrelated, and each diagonal element is the variance of a single component

Let  $X \sim N(0, R)$

$$E[\vec{x} \vec{x}^T] = R \text{ - by definition}$$

$$\hat{R} = \frac{1}{M} \sum_{k=1}^M \vec{x}_k \vec{x}_k^T$$

estimate  
(r.v.)

Law of large numbers  $\rightarrow \lim_{M \rightarrow \infty} \hat{R} \xrightarrow{\text{a.s.}} R$

, eigenvectors

$$R = E \Lambda E^T$$

\  
eigenvalues

$$E \text{ is orthonormal: } EE^T = E^T E = I$$

Gaussian distributions with zero mean have symmetry

$$E = [e_1, \dots, e_p]$$

$E$  can be viewed as a rigid body rotation

## Gaussian Random Variable Decorrelation

$$\tilde{\mathbf{x}} \sim N(\mathbf{0}, \mathbf{R}) ; \mathbf{R} = \mathbf{E} \Lambda \mathbf{E}^+$$

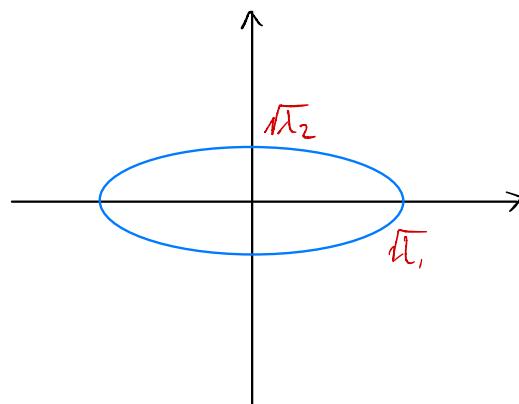
- Consider  $\tilde{\mathbf{x}} = \mathbf{E}^t \mathbf{x}$ , then

*transformed  
version of  $\mathbf{x}$*

*'rigid rotation'*

$$\begin{aligned}\mathbb{E}[\tilde{\mathbf{x}} \tilde{\mathbf{x}}^t] &= \mathbb{E}[\mathbf{E}^t \mathbf{x} \mathbf{x}^t \mathbf{E}] \\ &= \mathbf{E}^t \mathbb{E}[\mathbf{x} \mathbf{x}^t] \mathbf{E} \\ &= \mathbf{E}^t \mathbf{R} \mathbf{E} \quad \text{bring out E's as} \\ &\quad \text{they are deterministic} \\ &= \mathbf{E}^t \mathbf{E} \Lambda \mathbf{E}^t \mathbf{E} \\ &= \Lambda\end{aligned}$$

- Therefore, the elements of  $\tilde{\mathbf{x}}$  are uncorrelated with variance  $\mathbb{E}[\tilde{x}_i^2] = \Lambda_{ii}$ . , zero elsewhere  $\rightarrow$   $\tilde{\mathbf{x}}$  elements uncorrelated but not all same variance
- Therefore, • The elements of  $\tilde{\mathbf{x}}$  are independent, since  $\tilde{\mathbf{x}}$ , as the linear transform of  $\mathbf{x}$ , is Gaussian distributed.  $\rightarrow$  because  $\mathbf{x}$  is Gaussian & zero-mean



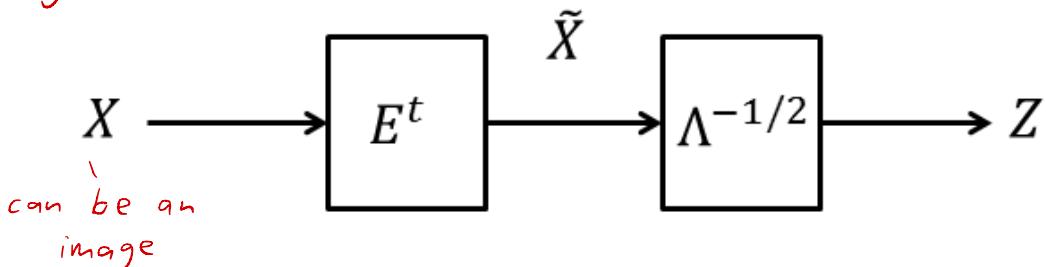
$$\begin{aligned}
 R &= E[xx^+] \\
 &= E\Lambda E^+ \\
 R \hat{=} \hat{R} &= \frac{1}{N} \sum_{i=1}^N y_i y_i^+ \quad \text{Sample covariance} \\
 \hat{R} &= \frac{1}{N} \tilde{X} \tilde{X}^+ = \frac{1}{N} \sum_{i=1}^N y_i y_i^+
 \end{aligned}$$

$\tilde{X} = \begin{bmatrix} y_1 & \dots & y_N \end{bmatrix}$  sample vectors  
 $10^6 \times 10^2$   
 $p \times N$   
each training image

Turns out that  
 $\tilde{X} = U \sum V^+$   
e-vecs of R  
e-vals of R  
 works for non-symmetric matrices

## Whitening Gaussian Random Variables

"Spherizing"

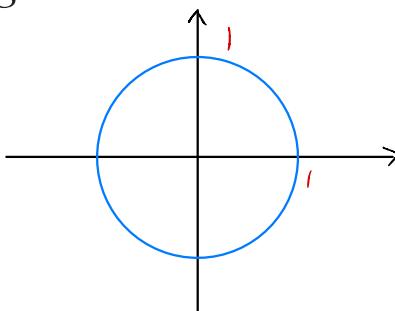


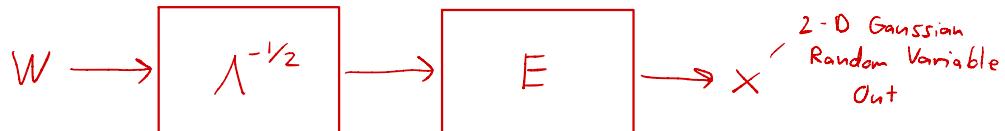
$$\mathbb{E}[Z Z] = I$$

- So  $E^t$  decorrelates  $x$ , while  $\Lambda^{-\frac{1}{2}}E^t$  whitens  $x$ .

$$E^t x = \begin{bmatrix} e_1^t \\ \vdots \\ e_p^t \end{bmatrix} \begin{bmatrix} x \end{bmatrix}$$

- The eigenvectors  $e_k$ , called eigen-signals, are basis vectors to represent the signal  $x$ .
- If  $x$  represents an image, then the eigenvectors  $e_k$  are also called *eigenimages*.





*Arbitrary Generator*

## Eigenimage Estimation

- Assume we have  $n$  training vectors  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , then

$$\begin{aligned} R_x &= E[\mathbf{x}_k \mathbf{x}_k^t] \\ &= E[XX^t]/n \\ &\cong XX^t/n \\ &= S \end{aligned}$$

where  $S = \frac{1}{n}XX^t$  is the sample correlation matrix

$$S_{ij} = \frac{1}{n} \sum_{k=1}^n X_{ik} X_{jk}$$

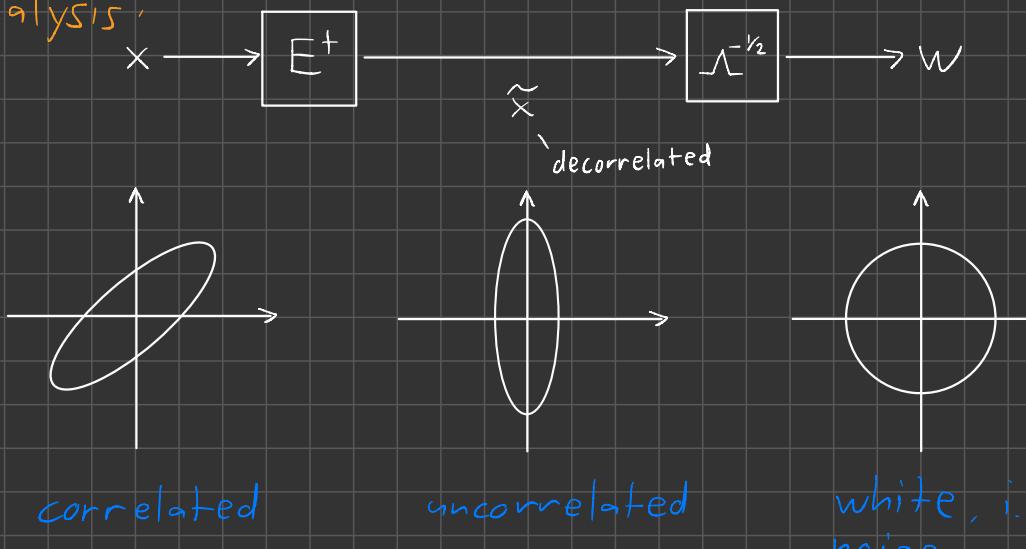
- Decompose  $S$  as

$$S = \hat{E} \hat{\Lambda} \hat{E}^t$$

where  $\hat{E}$  is an estimate of the eigenvectors, and  $\hat{\Lambda}$  is an estimate of the eigenvalues.

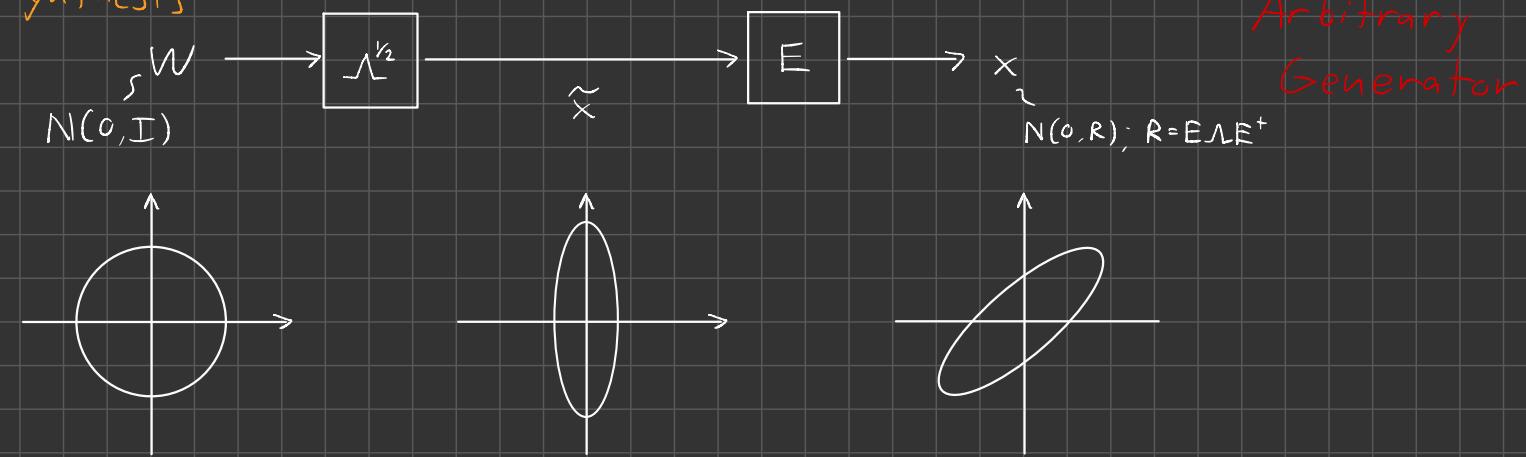
- $E$  could be very large, especially when  $X$  represents  $n$  images.

Analysis:



Hypothesis  
Testing

Synthesis:



Arbitrary  
Generator

$$X \sim N(0, R)$$

$$R = E \Lambda E^+$$

$$W = \Lambda^{-\frac{1}{2}} E^+ X ; \quad X = E \Lambda^{\frac{1}{2}} W$$

$$W \sim N(0, I)$$

$R = I$ , which is diagonal  $\rightarrow$  data points uncorrelated

$\hookrightarrow$  uncorrelated + Gaussian: data points independent

$W$  is a column vector of length  $p$ , just like  $X$

-generating  $W$  on a computer is straightforward

-components of  $\tilde{X}$  are 0-mean and independent, but they don't have the same variance/standard deviation

Ex: Generate  $X \sim N(\mu, R)$

Step 1)  $R = E \Lambda E^+$  - call routine to calculate  $E, \Lambda$  from  $R$

Step 2)  $W \sim (0, I)$  - call rand function  $p$  times to generate  $W$

Step 3)  $\tilde{X} = \Lambda^{1/2} W$      $\Lambda^{1/2}$  are the standard deviations?

Step 4)  $X' = \tilde{Y} = \sum_{i=1}^p e_i \tilde{x}_i$

Step 5)  $X = X' + \mu$

For large  $X \in \mathbb{R}$

## Singular Value Decomposition (SVD)

not the same as eigendecomposition

- For  $n < p$  it looks like

$$\begin{bmatrix} X \\ \cancel{X} \\ p \times n \end{bmatrix} = \begin{bmatrix} U \\ p \times n \end{bmatrix} \begin{bmatrix} \Sigma \\ n \times n \end{bmatrix} \begin{bmatrix} V^t \\ n \times n \end{bmatrix}$$

/ eigenimages are  
 the columns

raster-order  
 training images

- The columns of  $U$  are orthonormal and called left hand singular vectors.  $\rightarrow U^+U = I_n$
- The columns of  $V$  are orthonormal and called right hand singular vectors.  $\rightarrow V^+V = I_n$
- $\Sigma$  is diagonal matrix of singular values.

Any Linear Transformation

- 1) rigid body transformation #1
- 2) squeezing & stretching
- 3) rigid body transformation #2

$$\tilde{y} = \begin{bmatrix} U^+ \\ \vdots \end{bmatrix} \sim y \rightarrow \hat{y} = \begin{bmatrix} U \\ \vdots \end{bmatrix} \sim \tilde{y} = U\tilde{y}$$

$$\tilde{w} = T \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$E[\tilde{y}\tilde{y}^+] = R_y = \sum^2 \text{ like the Fourier coefficients}$$

$$\begin{aligned} E[\tilde{w}\tilde{w}^+] &= E[T w w^+ T^+] \\ &= T E[w w^+] T^+ \\ &= T I T^+ = I \end{aligned}$$

## Eigenimage Estimation using SVD

- Notice that

$$\begin{aligned} XX^t &= U\Sigma V^t V\Sigma U^t \\ &= U\Sigma^2 U^t \end{aligned}$$

So  $U$  is the set of the desired eigenvectors of  $XX^t$ .

- How to compute  $U$ ?

- Notice that  $X^t X = V\Sigma^2 V^t$  is a  $n \times n$  matrix, and  $V$  contains eigenvectors of a much smaller matrix.
- Algorithm

\* Find eigenvectors  $V$  of the matrix  $X^t X$

\* Compute  $XV = U\Sigma$ , Then the columns of  $U$  are the normalized columns of  $XV$ , and  $\Sigma$  are the normalization factors.

basis with decreasing weight of eigenvectors

\* Form  $B \leftarrow \frac{1}{N} X^t X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$

\*  $(V, \Sigma^2) \leftarrow \text{eigen}(B)$

\*  $\Sigma = \sqrt{\Sigma^2}$

\*  $U \leftarrow XVE^{-1}$  decorrelates normalizes

\* eigenvectors  $\leftarrow U$   
eigenvalues  $\leftarrow \Sigma^2$

represent  
the  
image

Common: Let  $x_1, \dots, x_n$  be i.i.d. R.V.

$X = [x_1, \dots, x_n]$  Experimentation:  $x_j$  could be an observation  
 $x_i$  are random vectors  $\sim N(0, R)$   
distribution that you sample from

$$E[x_i x_i^+] = E[\square] = R = E\left[\frac{1}{n} \sum_{i=1}^n x_i x_i^+\right]$$

Expectation is linear,  $R$  summed  $n$  times then scaled by  $\frac{1}{n}$  is simply  $R$

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^+ = \frac{1}{n} X X^+ = \hat{R} \text{ - estimate of } R \text{ (random)}$$

often:  $n < p$   $\Rightarrow X = \begin{bmatrix} \square \\ \vdots \\ \square \end{bmatrix}$  tall & thin  $R = E[\hat{R}]$   
 $\Rightarrow \hat{R}$  is an unbiased estimate

$$\hat{R} = \frac{1}{n} X X^+ = \begin{bmatrix} \square & \square & \cdots & \square \end{bmatrix} \text{ if } p = 10^6, \text{ then } \hat{R} \in \mathbb{R}^{p \times p}$$

$\hat{R}$  has  $10^{12}$  entries, not very practical

$$\hat{R} = U \Lambda U^+ ?? \rightarrow \text{enter SVD}$$

$$X = U \Sigma V^+$$

$$\begin{bmatrix} \square \\ \vdots \\ \square \end{bmatrix} = \begin{bmatrix} \square & \square & \square & \square \end{bmatrix} \begin{bmatrix} \square & \square & \square & \square \end{bmatrix}^T$$

- can always do this decomposition

$$\hat{R} = \frac{1}{n} X X^+ = U \Sigma^2 U^+ \rightarrow \text{the eigendecomposition of } \hat{R}$$

$$X^+ X = U \Sigma U^+ U \Sigma U^+ = U \Sigma^2 U^+$$

$\downarrow$   
smaller than  $X X^+$

Interpretation:  $X^+ X$  is taking the correlation of the images, and it is symmetric

## Overall Process Steps:

Step 1) Do SVD of  $X$   
 $X = U \Sigma V^+$ , can call in numpy

can effectively be considered  
training data

Step 2) Eigendecomposition  $\hat{R} = U \Sigma^2 V^+$

$$\boxed{U} \boxed{\Sigma^2} \boxed{V^+}$$

column space of  $U$   
does not span  $\mathbb{R}^p$

Step 3) Use them!

columns of  $V$  are Eigen vectors (Eigenimages in this case)

$$\tilde{X} = U^+ X$$

↑  
low-dimensional representation of image  $X$

$\left. \begin{array}{l} \text{image} \\ \text{dimensionality reduction technique} \end{array} \right\} \uparrow \boxed{\quad} = \boxed{\quad} \boxed{\quad}$

## Session 24 continuation from Session 23

$$X = U \boxed{\Sigma} V^+$$

↑  
singular values  
left-hand singular vectors  
(typically correlated) → training vectors

columns of  $X$  are the sample vectors

Correlated = not orthogonal

$$X = U \boxed{\Sigma}$$

V de-correlates the columns of  $X$ , resulting columns will be orthogonal but won't have the same norm (some columns will have more energy than others)

$$X = U \boxed{\Sigma}$$

Multiply by diagonal matrix → scale each column  
If multiplication on LHS → scale each row (as FYI)  
 $\Sigma'$  scales the de-correlated columns of  $X$  so they all have the same unit norm

Intuitively, this is sort of like Gram-Schmidt

Remember:  $X X^T / X^T X$  is exactly the correlation between the rows / columns of  $X$ , respectively

$$X = U \Sigma V^+$$

$$\hookrightarrow X^+ X = V \Sigma^2 V^+ = R$$

"small" matrix,  $n \times n$

this is how  $V$  can be calculated

This is how  
 $V$  is calculated



then, decorrelate the columns of  $X$  \*

then, scale with  $\Sigma^{-1}$  to get  $V$

-columns of  $V$  are the eigenimages

-if singular values are ordered from largest to smallest, which is pretty conventional, you'll get the eigenimages in order from lowest frequency to highest frequency - why is this??

### Algorithm:

Step 0) Form matrix of training data, used to build model

$X = [x, \dots, x_n]$  → could be images, signals, etc.

Step 1) Compute  $V$

a)  $R \leftarrow X^+ X$

b)  $R = V \Sigma^2 V^+$

c)  $V = X V \Sigma^{-1}$

Step 2) Use the eigenimages!!

$$U = [u_1, \dots, u_n]$$

Take new image  $y$

$$f_i = u_i^T y = \langle u_i, y \rangle$$

scalar feature

$$f = U^+ y \rightsquigarrow \begin{bmatrix} & \\ x_1' & n \times p \\ & \end{bmatrix} = \begin{bmatrix} & \\ & \\ & \end{bmatrix} \begin{bmatrix} & \\ p \times 1 & \end{bmatrix}$$

intuition:  $f$  is an  $n$ -dimensional representation of a  $p$ -dimensional image

dimensionality reduction process

$$\tilde{y} = Uf = y + \varepsilon \quad \text{transforming back}$$

$$\begin{bmatrix} \tilde{y} \\ \vdots \\ \tilde{y} \end{bmatrix} = \begin{bmatrix} U & f \end{bmatrix} \begin{bmatrix} y \\ \vdots \\ y \end{bmatrix} + \varepsilon$$

f is a set of scalar weights that multiply each column of U, which are the basis vectors, and the result is a representation of y with some residual error

$p \times 1$   $p \times n$

Intuitively sort of like a Fourier Transform, or an image compression algo!